# Demonstrating Hybrid Provenance Queries from Script Runs

Yang Cao[1], Peter Slaughter[5], Christopher Jones[5], Matthew B. Jones[5], Qiwen Wang[2], Duc Vu[3],
Priyaa Thavasimani[4], Qian Zhang[1], Timothy McPhillips[1], Paolo Missier[4], Lauren Walker[5],
Dave Vieglais[6], Bertram Ludäscher[1 2]

[1]School of Information Sciences, University of Illinois at Urbana-Champaign
[2]Department of Computer Science, University of Illinois at Urbana-Champaign
[3]Department of Electrical and Computer Engineering, University of Illinois at Chicago
[4]School of Computing Science, Newcastle University, UK
[5]National Center for Ecological Analysis and Synthesis, UCSB, Santa Barbara, USA
[6] University of Kansas, Lawrence, USA

*Data provenance* is metadata that describes the origin and processing history of a data artifact. In the computational and data sciences, we can distinguish two forms of data provenance, i.e., *prospective* and *retrospective* provenance. The former describes the general *workflow* by which data is produced, while the latter consists of runtime observables, e.g., names, locations, and contents of files read and written during a script run (coarse-grained observables). Fine-grained provenance can be captured as well, e.g., changes to individual data records, program variables, parameter settings, etc. By *hybrid provenance,* we mean a combination of both provenance types, i.e., detailed retrospective provenance information is "situated" within the context of the conceptual-level workflow (prospective provenance), thus aiming to get the "best of both worlds": The appendix depicts such a hybrid graph with alternating process steps and data elements. The latter can be instantiated to document the specific files used in a concrete workflow run (see Appendix). Since computational and data science experiments can often last days, weeks, or even months and often require the execution of multiple scripts or workflows with varying input datasets and parameters, some of these script runs appear as chained together implicitly via intermediate data files, i.e., the output of one script run was subsequently used as an input to another script run. In this way, larger ad-hoc dataflow graphs are generated as part of a data-driven, exploratory workflow. By *multi-run provenance* we mean such dataflow graphs that link individual script runs via intermedia data files.

Our system demonstration (extending earlier work [1]; see also [2]) will illustrate the variety of provenance information that we are able to capture, query, and visualize using a combination of tools for exposing both prospective and retrospective provenance. We show how prospective provenance can be declared using YesWorkflow (YW) annotations that reveal the fine-grained (variable level) dataflow graph implicit in scripts, and how this prospective provenance can be integrated with the coarse-grained (file-level) retrospective provenance information recorded by the DataONE Run Managers for MATLAB and R. We demonstrate the usefulness of integrating prospective and retrospective provenance in this way with queries:

1. **Prospective provenance queries in the context of a single script.** This can expose and test data dependencies at the workflow-level.
2. **Retrospective provenance queries in the context of a single run of a single script:** captures actual input and output files of a script run and other runtime observables.
3. **Hybrid provenance query in the context of a single script and single run:** blends retrospective and prospective provenance, yielding new knowledge artefacts.
4. **Provenance query in the context of multiple scripts and multiple runs:** query and visualize data dependencies across multiple script runs

Our demonstration queries and provenance reports thus yield a more complete and comprehensible picture of data provenance from multiple script runs.

## References

[1] Y Cao, D Vu, Q Wang, Q Zhang, P Thavasimani, T McPhillips, P Missier, B Ludäscher. YesWorkflow. Demonstration repository. https://github.com/idaks/dataone-ahm-2016-poster

[2] Q Zhang, Y Cao, Q Wang, D Vu, P Thavasimani, T McPhillips, P Missier, B Ludäscher. Revealing the Detailed Lineage of Script Outputs Using Hybrid Provenance. Submitted to IDCC 2017 (Practice Paper track).

**Appendix 1: Ocean Health Index for Howe Sound, British Columbia Example**

1. Prospective provenance graphs produced by YW (see Figure 1)
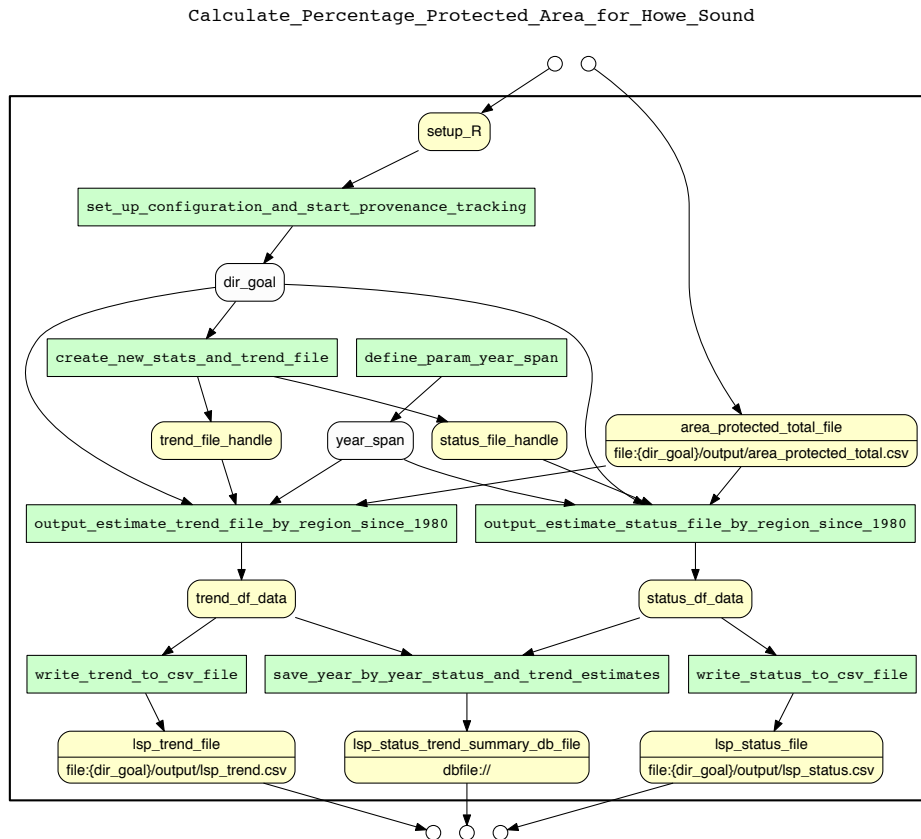2. Hybrid provenance graphs produced by YW and DataONE RunManager R client (see Figure 2)



**Figure** 1 Complete YW prospective provenance graph for the `OHIBC_Howe_Sound` example. The diagram shows the workflow for calculating the "lasting special places sub goal" of the Ocean Health Index for Howe Sound, British Columbia that is based on the percentage of protected area within 3km of shore and 1km buffer zone inland, and how that has changed over time. The green box represents a computation step and a round yellow square box represents a data element declared in the script via YW tags.
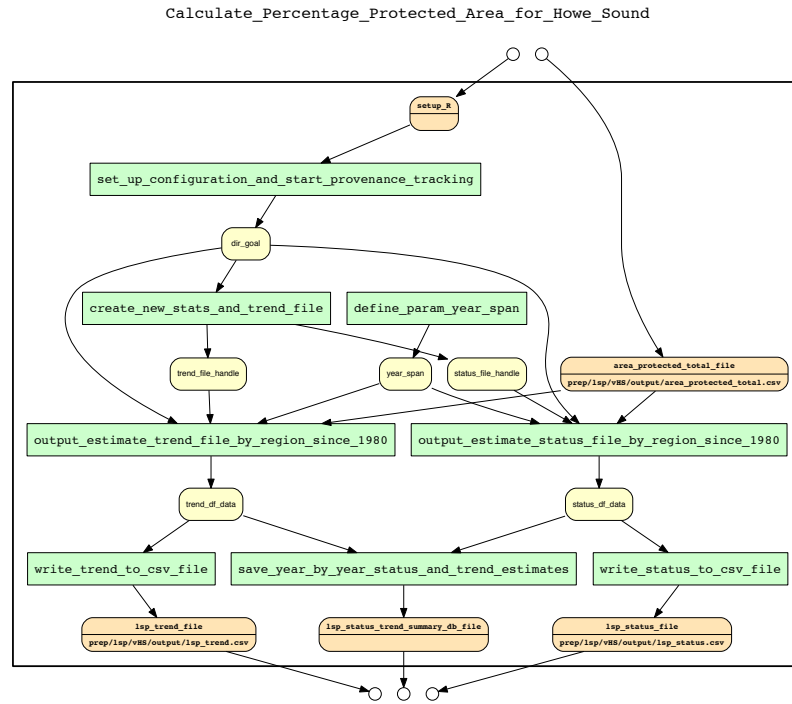
**Figure2** Complete reconstructed hybrid provenance graph for the `OHIBC_Howe_Sound` example using YW and DataONE RunManager R client. In the diagram, the runtime file-level observables are augmented with YW prospective graph.

**Appendix 2: Carbon 3/Carbon 4 (C3/C4) soil maps for North America example**

1. Hybrid provenance graph produced by YW and RunManager MATLAB client (see Figure 3)
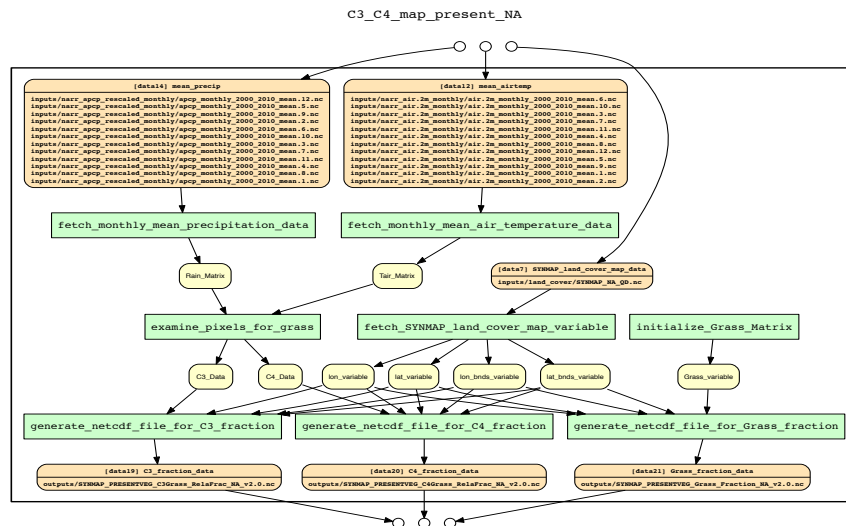2. Retrospective provenance query in the context of multiple scripts and multiple runs (see Figure 4)



**Figure** 3 Complete reconstructed hybrid provenance graph for the `C3C4` example (version 1) using YW and DataONE RunManager MATLAB client. The example script (in MATLAB) produces Carbon 3/Carbon 4 (C3/C4) soil maps for North America using average rain and air temperature monthly data from year 2000 to 2010.
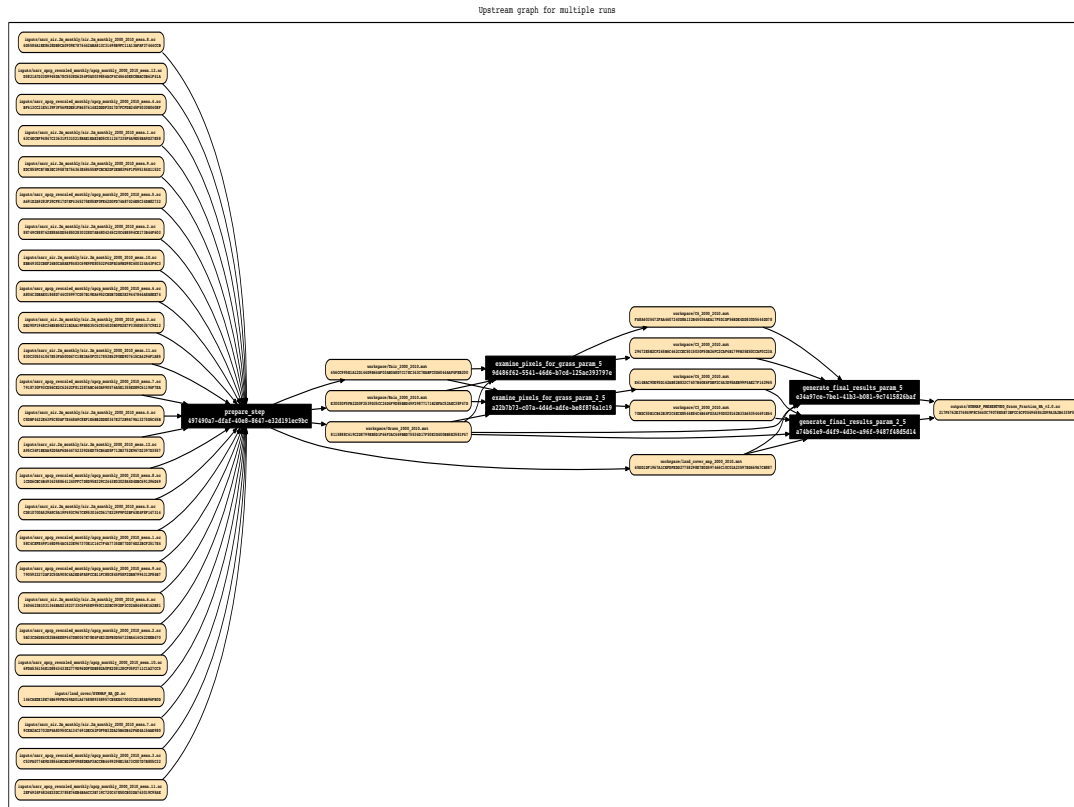
**Figure** 4 Upstream retrospective provenance subgraph (longitudinal view) recording multiple scripts and multiple runs for `C3C4` example. The generated graph recorded a set of runs that consist of three scripts, five runs and the second script was run with two different parameter values (i.e., 2.5 and 5) twice. The example script set is called `C3C4` example version 2. In the diagram, an orange yellow box represents a data file and a black box represents an execution (a run). When two data files sharing with a single orange yellow box, it indicates that these two files have the identical contents.